

GLOBAL CONVERGENCE OF RADIAL BASIS FUNCTION TRUST-REGION ALGORITHMS FOR DERIVATIVE-FREE OPTIMIZATION*

STEFAN M. WILD[†] AND CHRISTINE SHOEMAKER[‡]

Abstract. We analyze globally convergent, derivative-free trust-region algorithms relying on radial basis function interpolation models. Our results extend the recent work of Conn, Scheinberg, and Vicente [*SIAM J. Optim.*, 20 (2009), pp. 387–415] to fully linear models that have a nonlinear term. We characterize the types of radial basis functions that fit in our analysis and thus show global convergence to first-order critical points for the ORBIT algorithm of Wild, Regis, and Shoemaker [*SIAM J. Sci. Comput.*, 30 (2008), pp. 3197–3219]. Using ORBIT, we present numerical results for different types of radial basis functions on a series of test problems. We also demonstrate the use of ORBIT in finding local minima on a computationally expensive environmental engineering problem involving remediation of contaminated groundwater.

Key words. derivative-free optimization, radial basis functions, trust-region methods, nonlinear optimization

AMS subject classifications. 65D05, 90C30, 90C56

DOI. 10.1137/

1. Introduction. This paper concerns algorithms for solving the unconstrained optimization problem

$$(1.1) \quad \min_{x \in \mathbb{R}^n} f(x),$$

when only function values $f(x)$ (and not values of the gradient $\nabla f(x)$ or higher-order derivatives) are available to the optimization algorithm. This situation typically arises when evaluation of the function f requires running a numerical simulation or performing a physical experiment. It is important to distinguish the present setting from *nonsmooth optimization*, which concerns problems where derivatives of f do not exist. Here we consider deterministic, real-valued functions f that are assumed to be continuously differentiable with a Lipschitz gradient ∇f and bounded from below.

Apart from the zero-order (“derivative-free”) methods examined here, there are two primary approaches for solving (1.1) when derivatives of f are not directly available. When the source code for evaluating f is available, algorithmic differentiation (AD) [15] can be used to produce a derivative code, which can then be used in derivative-based optimization methods. Similarly, numerical differentiation (e.g.,

*Published electronically May 8, 2013. This paper is a modified version of “Global Convergence of Radial Basis Function Trust Region Derivative-Free Algorithms,” which originally appeared in *SIAM Journal on Optimization*, Volume 21, Number 3, 2011, pages 761–781. This paper has been created by UChicago Argonne, LLC, Operator of Argonne National Laboratory (“Argonne”). Argonne, a U.S. Department of Energy Office of Science laboratory, is operated under Contract No. DE-AC02-06CH11357.

<http://www.siam.org/journals/sirev/55-2/xxxxx.html>

[†]Mathematics and Computer Science Division, Argonne National Laboratory, Argonne, IL 60439 (wild@mcs.anl.gov). The work of this author was supported by a U.S. Department of Energy Computational Science Graduate Fellowship under grant DE-FG02-97ER25308.

[‡]School of Civil and Environmental Engineering and School of Operations Research and Information Engineering, Cornell University, Hollister Hall, Ithaca, NY 14853 (cas12@cornell.edu). The work of this author was supported by NSF grants BES-022917, CBET-0756575, CCF-0305583, and DMS-0434390.

with a finite-difference scheme) can be employed to obtain derivative approximations for use in derivative-based methods. The former approach is not viable, for example, when some portion of the function evaluation corresponds to a black box (e.g., requires a proprietary/executable-only evaluation). The latter approach can require careful selection of finite-difference parameters (see, e.g., [27]), and performing $\mathcal{O}(n)$ function evaluations at every iteration in order to determine a full gradient may be prohibitively expensive.

Research in *derivative-free optimization* has received renewed and sustained interest over the past 15 years (see the recent books [10, 20] for a summary). The proliferation of high-performance computing and an abundance of legacy codes have driven demand for derivative-free methods throughout computational science and engineering. These methods can generally be categorized into those based on systematic sampling of the function along well-chosen directions [1, 18, 21, 22, 23], and those employing a trust-region framework with a local approximation of the function [7, 24, 30, 31, 32].

Many methods in the former category are popular with engineers for the ease of basic implementation. These include the Nelder-Mead simplex algorithm [23] and pattern search [22]. Such methods also admit natural parallel implementations [18, 21], where different poll directions are sent to different processors for evaluation, and extensions can be shown to converge even when derivatives do not exist [1].

Methods in the latter category (including the one analyzed in this paper) use prior function evaluations to construct a model that approximates the function in a neighborhood of the current iterate. These models (for example, fully quadratic [7, 24, 31], underdetermined or structured quadratic [32], or radial basis functions (RBFs) [30, 38]) yield computationally attractive derivatives and hence are easy to optimize over within the neighborhood.

Each of these methods, as well as the one analyzed in this paper, assumes that the function f is a black box. Knowledge of additional structure in the function f can be exploited in order to reduce the number of evaluations of f required to solve (1.1). Recent examples of exploiting known structure include methods for nonlinear least-squares problems with black-box residuals [28] and methods for cases where the unavailable derivatives are sparse [2].

A keystone of the present work is our assumption that the computational expense of the function evaluation yields a bottleneck for optimization (the expense of evaluating the function at a single point outweighing any other expense or overhead of an algorithm). In some applications this could mean that function evaluation can require from a few seconds on a state-of-the-art machine to several hours on a large cluster, even when the functions are parallelized. The functions that motivate our work usually depend on complex deterministic computer simulations, including those that numerically solve systems of PDEs governing underlying physical phenomena.

This paper is driven by work on the ORBIT (Optimization by Radial Basis functions In Trust regions) algorithm [38] and provides the key theoretical conditions needed for such algorithms to converge to first-order critical points. We find that the popular thin-plate spline RBFs do not fit in this globally convergent framework. Furthermore, our numerical results comparing RBF types show that the Gaussian RBFs popularly used in kriging [11, 19] are not as effective in our algorithms as are other RBF types. Comparisons with other types of derivative-free methods can be found in [38].

ORBIT is a trust-region algorithm relying on an interpolating RBF model with

a *linear polynomial tail*. A primary distinction between ORBIT and the previously proposed RBF-based algorithm in [30] is the management of the interpolation set (Algorithm 3). In contrast to [30], the expense of our objective function allows us to effectively ignore the computational complexity of the overhead of building and maintaining the RBF model.

Our first goal is to show global convergence to first-order critical points for very general interpolation models. In section 2 we review the multivariate interpolation problem and show that the local error between the function (and its gradient) and an interpolation model (and its gradient) can be bounded by using a simple condition on $n + 1$ of the interpolation points. In the spirit of [8], we refer to such interpolation models as *fully linear*. In section 3 we review derivative-free trust-region methods and analyze conditions necessary for global convergence when fully linear models are employed. For this convergence analysis we benefit from the results in [9].

Our next goal is to use this analysis to identify the conditions necessary for obtaining a globally convergent trust-region method by using an interpolating RBF-based model. In section 4 we introduce radial basis functions and the fundamental property of *conditional positive definiteness*, which we rely on in ORBIT to construct uniquely defined RBF models with bounded coefficients. We also give necessary and sufficient conditions for different RBF types to fit within our framework.

In section 5 we examine the effect of selecting from three popular radial basis functions covered by the theory by running the resulting algorithm on a set of smooth test functions. We also examine the effect of varying the maximum number of interpolation points. We motivate the use of ORBIT to quickly find local minima of computationally expensive functions with an application problem (requiring nearly 1 CPU-hour per evaluation on a Pentium 4 machine) arising from detoxification of contaminated groundwater. We note that additional computational results, both on a set of test problems and on two applications from environmental engineering, as well as more practical considerations, are addressed in [38].

2. Interpolation Models. We begin our discussion on models that interpolate a set of scattered data with an introduction to the polynomial models that are heavily utilized by derivative-free trust-region methods in the literature [7, 24, 31, 32].

2.1. Notation. We first collect the notation conventions used throughout the paper. \mathbb{N}_0^n will denote n -tuples from the natural numbers including zero. A vector $x \in \mathbb{R}^n$ will be written in component form as $x = [\chi_1, \dots, \chi_n]^T$ to differentiate it from a particular point $x_i \in \mathbb{R}^n$. For $d \in \mathbb{N}_0$, let \mathcal{P}_{d-1}^n denote the space of n -variate polynomials of total degree no more than $d - 1$, with the convention that $\mathcal{P}_{-1}^n = \emptyset$. Let $\mathcal{Y} = \{y_1, y_2, \dots, y_{|\mathcal{Y}|}\} \subset \mathbb{R}^n$ denote an interpolation set of $|\mathcal{Y}|$ points where (y_i, f_i) is known. For ease of notation, we will often assume interpolation relative to some base point $x_b \in \mathbb{R}^n$, made clear from the context, and will employ the set notation $x_b + \mathcal{Y} = \{x_b + y : y \in \mathcal{Y}\}$. We will work with a general norm $\|\cdot\|_k$ that we relate to the 2-norm $\|\cdot\|$ through a constant c_1 , depending only on n , satisfying

$$(2.1) \quad \|\cdot\| \leq c_1 \|\cdot\|_k, \quad \forall k.$$

The polynomial interpolation problem is to find a polynomial $P \in \mathcal{P}_{d-1}^n$ such that

$$(2.2) \quad P(y_i) = f_i, \quad \forall y_i \in \mathcal{Y},$$

for arbitrary values $f_1, \dots, f_{|\mathcal{Y}|} \in \mathbb{R}$. Spaces where unique polynomial interpolation is always possible given an appropriate number of distinct data points are called *Haar*

spaces. A classic theorem of Mairhuber and Curtis (see [35, p. 19]) states that Haar spaces do not exist when $n \geq 2$. Hence additional conditions are necessary for the multivariate problem (2.2) to be well-posed. We use the following definition.

DEFINITION 2.1. *The points \mathcal{Y} are \mathcal{P}_{d-1}^n -unisolvent if the only polynomial in \mathcal{P}_{d-1}^n that vanishes at all points in \mathcal{Y} is the zero polynomial.*

The monomials $\{\chi_1^{\alpha_1} \cdots \chi_n^{\alpha_n} : \alpha \in \mathbb{N}_0^n, \sum_{i=1}^n \alpha_i \leq d-1\}$ form a basis for \mathcal{P}_{d-1}^n , and hence any polynomial $P \in \mathcal{P}_{d-1}^n$ can be written as a linear combination of such monomials. In general, for a basis $\pi : \mathbb{R}^n \rightarrow \mathbb{R}^{\hat{p}}$ we will use the representation $P(x) = \sum_{i=1}^{\hat{p}} \nu_i \pi_i(x)$, where $\hat{p} = \dim \mathcal{P}_{d-1}^n = \binom{n+d-1}{n}$. Hence finding an interpolating polynomial $P \in \mathcal{P}_{d-1}^n$ is equivalent to finding coefficients $\nu \in \mathbb{R}^{\hat{p}}$ for which (2.2) holds.

Defining $\Pi \in \mathbb{R}^{\hat{p} \times |\mathcal{Y}|}$ by $\Pi_{i,j} = \pi_i(y_j)$, it follows that \mathcal{Y} is \mathcal{P}_{d-1}^n -unisolvent if and only if Π is full rank, $\text{rank} \Pi = \hat{p}$. Furthermore, the interpolation in (2.2) is unique for arbitrary right-hand-side values $f_1, \dots, f_{|\mathcal{Y}|} \in \mathbb{R}$ if and only if $|\mathcal{Y}| = \hat{p}$ and Π is nonsingular. In this case, the unique polynomial is defined by the coefficients $\nu = \Pi^{-T} f$.

One can easily see that existence and uniqueness of an interpolant are independent of the particular basis π employed. However, the conditioning of the corresponding matrix Π depends strongly on the basis chosen, as noted (for example) in [8].

Based on these observations, we see that in order to uniquely fit a polynomial of degree $d-1$ to a function, at least $\hat{p} = \dim \mathcal{P}_{d-1}^n = \binom{n+d-1}{n}$ function values must be known. When n is not very small, the computational expense of evaluating f to repeatedly fit even a quadratic (with $\hat{p} = \frac{(n+1)(n+2)}{2}$) is large.

2.2. Fully Linear Models. We now explore a class of *fully linear* interpolation models, which can be formed by using as few as a linear (in the dimension, n) number of function values. Since such models are heavily tied to Taylor-like error bounds, we will require assumptions on the function f as in this definition from [8].

DEFINITION 2.2. *Suppose that $\mathcal{B} = \{x \in \mathbb{R}^n : \|x - x_b\|_k \leq \Delta\}$ and $f \in C^1[\mathcal{B}]$. For fixed $\kappa_f, \kappa_g > 0$, a model $m \in C^1[\mathcal{B}]$ is said to be fully linear on \mathcal{B} if for all $x \in \mathcal{B}$*

$$(2.3) \quad |f(x) - m(x)| \leq \kappa_f \Delta^2,$$

$$(2.4) \quad \|\nabla f(x) - \nabla m(x)\| \leq \kappa_g \Delta.$$

This definition ensures that first-order Taylor-like bounds exist for the model within the compact neighborhood \mathcal{B} . For example, if $f \in C^1[\mathbb{R}]$, ∇f has Lipschitz constant γ_f and if m is the derivative-based linear model $m(x_b + s) = f(x_b) + \nabla f(x_b)^T s$, then m is fully linear with constants $\kappa_g = \kappa_f = \gamma_f$ on any bounded region \mathcal{B} .

Since the function's gradient is unavailable in our setting, our focus is on models that interpolate the function at a set of points:

$$(2.5) \quad m(x_b + y_i) = f(x_b + y_i) \quad \text{for all } y_i \in \mathcal{Y} = \{y_1 = 0, y_2, \dots, y_{|\mathcal{Y}|}\} \subset \mathbb{R}^n.$$

Although we may have interpolation at more points, for the moment we work with a subset of exactly $n+1$ points and always enforce interpolation at the base point x_b so that $y_1 = 0 \in \mathcal{Y}$. The remaining n (nonzero) points compose the square matrix $Y = \begin{bmatrix} y_2 & \cdots & y_{n+1} \end{bmatrix}$.

We can now state error bounds, similar to those in [8], for our models of interest.

THEOREM 2.3. *Suppose that f and m are continuously differentiable in $\mathcal{B} = \{x : \|x - x_b\|_k \leq \Delta\}$ and that ∇f and ∇m are Lipschitz continuous in \mathcal{B} with Lipschitz constants γ_f and γ_m , respectively. Further suppose that m satisfies the interpolation*

conditions in (2.5) at a set of points $\{y_1 = 0, y_2, \dots, y_{n+1}\} \subseteq \mathcal{B} - x_b$ such that $\|Y^{-1}\| \leq \frac{\Lambda_Y}{c_1 \Delta}$, for a fixed constant $\Lambda_Y < \infty$ and c_1 from (2.1). Then for any $x \in \mathcal{B}$,

$$(2.6) \quad |f(x) - m(x)| \leq \sqrt{n} c_1^2 (\gamma_f + \gamma_m) \left(\frac{5}{2} \Lambda_Y + \frac{1}{2} \right) \Delta^2 = \kappa_f \Delta^2,$$

$$(2.7) \quad \|\nabla f(x) - \nabla m(x)\| \leq \frac{5}{2} \sqrt{n} \Lambda_Y c_1 (\gamma_f + \gamma_m) \Delta = \kappa_g \Delta.$$

Proved in [36], Theorem 2.3 provides the constants $\kappa_f, \kappa_g > 0$ such that conditions (2.3) and (2.4) are satisfied, and hence m is fully linear in a neighborhood \mathcal{B} containing the $n + 1$ interpolation points. This result holds for very general interpolation models, requiring only a minor degree of smoothness and conditions on the points being interpolated. The conditions on the interpolation points are equivalent to requiring that the points $\{y_1, y_2, \dots, y_{n+1}\}$ are *sufficiently affinely independent* (or equivalently, that the set $\{y_2 - y_1, \dots, y_{n+1} - y_1\}$ is sufficiently linearly independent), with Λ_Y quantifying the degree of independence.

One can iteratively construct a set of such points given a set of candidate displacements $\mathcal{D} = \{d_1, \dots, d_{|\mathcal{D}|}\} \subset \{x \in \mathbb{R}^n : \|x\|_k \leq \Delta\}$ (e.g., from the nearby points $x_b + \mathcal{D}$ at which f has been evaluated) by using LU- and QR-like algorithms as noted in [8].

For example, in ORBIT points are added to the interpolation set \mathcal{Y} one at a time by using a QR-like variant described in [38]. The crux of the algorithm is to add a candidate from \mathcal{D} to \mathcal{Y} if its projection onto the subspace orthogonal to $\text{span} \mathcal{Y}$ is sufficiently large (as measured by a constant $\theta \in (0, 1]$). If the candidates in \mathcal{D} are not sufficiently affinely independent, such algorithms also produce points belonging to \mathcal{B} that are perfectly conditioned with respect to the projection so that m can be easily made fully linear in fewer than n function evaluations.

We conclude this section by stating a lemma from [38] that ensures a QR-like procedure similar to one mentioned yields a set of points in \mathcal{Y} satisfying $\|Y^{-1}\| \leq \frac{\Lambda_Y}{c_1 \Delta}$.

LEMMA 2.4. *Let $QR = \frac{1}{c_1 \Delta} Y$ denote a QR factorization of a matrix $\frac{1}{c_1 \Delta} Y$ whose columns satisfy $\|\frac{Y_j}{c_1 \Delta}\| \leq 1$, $j = 1, \dots, n$. If $r_{ii} \geq \theta > 0$ for $i = 1, \dots, n$, then $\|Y^{-1}\| \leq \frac{\Lambda_Y}{c_1 \Delta}$ for a constant Λ_Y depending only on n and θ .*

3. Derivative-Free Trust-Region Methods. The interpolation models of the previous section were constructed to approximate a function in a local neighborhood of a point x_b . The natural algorithmic extensions of such models are trust-region methods (given full treatment in [6]), whose general form we now briefly review.

Trust-region methods generate a sequence of iterates $\{x_k\}_{k \geq 0} \subseteq \mathbb{R}^n$ by employing a surrogate model $m_k : \mathbb{R}^n \rightarrow \mathbb{R}$, assumed to approximate f within a neighborhood of the current x_k . For a (center, radius) pair $(x_k, \Delta_k > 0)$ we define the *trust region*

$$(3.1) \quad \mathcal{B}_k = \{x \in \mathbb{R}^n : \|x - x_k\|_k \leq \Delta_k\},$$

where we distinguish the trust-region norm (at iteration k), $\|\cdot\|_k$, from other norms used here. New points are obtained by solving subproblems of the form

$$(3.2) \quad \min_s \{m_k(x_k + s) : x_k + s \in \mathcal{B}_k\}.$$

The pair (x_k, Δ_k) is then updated according to the ratio of actual to predicted improvement. Given a maximum radius Δ_{\max} , the design of the trust-region algorithm

ensures that f is sampled only within the relaxed level set

$$(3.3) \quad \mathcal{L}(x_0) = \{y \in \mathbb{R}^n : \|x - y\|_k \leq \Delta_{\max} \text{ for some } x \text{ with } f(x) \leq f(x_0)\}.$$

Hence one really requires only that f be sufficiently smooth within $\mathcal{L}(x_0)$.

When derivatives are unavailable, smoothness of the function f is no longer sufficient for guaranteeing that a model m_k approximates the function locally. Hence the main difference between classical and derivative-free trust-region algorithms is the addition of safeguards to account for and improve models of poor quality.

Historically (see [7, 24, 31, 32]), the most frequently used model is a quadratic,

$$(3.4) \quad m_k(x_k + s) = f(x_k) + g_k^T s + \frac{1}{2} s^T H_k s,$$

the coefficients g_k and H_k being found by enforcing interpolation as in (2.5). As discussed in section 2, these models rely heavily on results from multivariate interpolation. Quadratic models are attractive in practice because the resulting subproblem in (3.2), for a 2-norm trust region, is one of the only nonlinear programs for which *global* solutions can be efficiently computed [25].

A downside of quadratic models in our computationally expensive setting is that the number of interpolation points (and hence function evaluations) required is quadratic in the dimension of the problem. Noting that it may be more efficient to use function evaluations for forming subsequent models, Powell designed his **NEWUOA** code [32] to rely on least-change quadratic models interpolating fewer than $\frac{(n+1)(n+2)}{2}$ points. More recent work in [12, 14] has also explored loosening the restrictions of a quadratic number of geometry conditions.

3.1. Fully Linear Derivative-Free Models. Recognizing the difficulty (and possible inefficiency) of maintaining geometric conditions on a quadratic number of points, we will focus on using the fully linear models introduced in section 2. These models can be formed with a linear number of points while still maintaining the local approximation bounds in (2.3) and (2.4).

We will follow the general trust-region algorithmic framework introduced for linear models by Conn et al. [9] in order to arrive at a similar convergence result for the types of models considered here. Given standard trust-region inputs $0 \leq \eta_0 < \eta_1 < 1$, $0 < \gamma_0 < 1 < \gamma_1$, $0 < \Delta_0 \leq \Delta_{\max}$, and $x_0 \in \mathbb{R}^n$ and constants $\kappa_d \in (0, 1)$, $\kappa_f > 0$, $\kappa_g > 0$, $\epsilon > 0$, $\mu > \beta > 0$, $\alpha \in (0, 1)$, the general first-order derivative-free trust-region algorithm is shown in Algorithm 1. This algorithm is discussed in [9], and we note that it forms an infinite loop, a recognition that termination in practice is a result of exhausting a budget of expensive function evaluations.

A benefit of working with more general fully linear models is that they allow for nonlinear modeling of f . Hence, we will be interested primarily in models with nontrivial Hessians, $\nabla^2 m_k \neq 0$, which are uniformly bounded by some constant κ_H .

The sufficient decrease condition that we will use in Step 1.2 then takes the form

$$(3.7) \quad m_k(x_k) - m_k(x_k + s) \geq \frac{\kappa_d}{2} \|\nabla m_k(x_k)\| \min \left\{ \frac{\|\nabla m_k(x_k)\|}{\kappa_H}, \frac{\|\nabla m_k(x_k)\|}{\|\nabla m_k(x_k)\|_k} \Delta_k \right\},$$

for some prespecified constant $\kappa_d \in (0, 1)$. This condition is similar to those found in the trust-region setting when general norms are employed [6]. The following lemma guarantees we will always be able to find an approximate solution, s_k , to the subproblem (3.2) that satisfies condition (3.7).

Algorithm 1 Iteration k of a first-order (fully linear) derivative-free algorithm [9].

1.1. Criticality test If $\|\nabla m_k(x_k)\| \leq \epsilon$ and either m_k is not fully linear in \mathcal{B}_k or $\Delta_k > \mu \|\nabla m_k(x_k)\|$:

Set $\tilde{\Delta}_k = \Delta_k$ and make m_k fully linear on $\{x : \|x - x_k\|_k \leq \tilde{\Delta}_k\}$.

While $\tilde{\Delta}_k > \mu \|\nabla m_k(x_k)\|$:

Set $\tilde{\Delta}_k \leftarrow \alpha \tilde{\Delta}_k$ and make m_k fully linear on $\{x : \|x - x_k\|_k \leq \tilde{\Delta}_k\}$.

Update $\Delta_k = \max\{\tilde{\Delta}_k, \beta \|\nabla m_k(x_k)\|\}$.

1.2. Obtain trust-region step s_k satisfying a sufficient decrease condition (e.g., (3.7)).

1.3. Evaluate $f(x_k + s_k)$.

1.4. Adjust trust region according to ratio $\rho_k = \frac{f(x_k) - f(x_k + s_k)}{m_k(x_k) - m_k(x_k + s_k)}$:

$$(3.5) \quad \Delta_{k+1} = \begin{cases} \min\{\gamma_1 \Delta_k, \Delta_{\max}\} & \text{if } \rho_k \geq \eta_1 \text{ and } \Delta_k < \beta \|\nabla m_k(x_k)\| \\ \Delta_k & \text{if } \rho_k \geq \eta_1 \text{ and } \Delta_k \geq \beta \|\nabla m_k(x_k)\| \\ \Delta_k & \text{if } \rho_k < \eta_1 \text{ and } m_k \text{ is not fully linear} \\ \gamma_0 \Delta_k & \text{if } \rho_k < \eta_1 \text{ and } m_k \text{ is fully linear,} \end{cases}$$

$$(3.6) \quad x_{k+1} = \begin{cases} x_k + s_k & \text{if } \rho_k \geq \eta_1 \\ x_k + s_k & \text{if } \rho_k > \eta_0 \text{ and } m_k \text{ is fully linear} \\ x_k & \text{else.} \end{cases}$$

1.5. Improve m_k if $\rho_k < \eta_1$ and m_k is not fully linear.

1.6. Form new model m_{k+1} .

LEMMA 3.1. If $m_k \in C^2(\mathcal{B}_k)$ and $\kappa_H > 0$ satisfies

$$(3.8) \quad \infty > \kappa_H \geq \max_{x \in \mathcal{B}_k} \|\nabla^2 m_k(x)\|,$$

then for any $\kappa_d \in (0, 1)$ there exists an $s \in \mathcal{B}_k - x_k$ satisfying (3.7).

Lemma 3.1 (proved in [36]) is our variant of similar ones in [6] and describes a back-tracking line search algorithm to obtain a step that yields a model reduction at least a fraction of that achieved by the Cauchy point. As an immediate corollary we have that there exists a step $s \in \mathcal{B}_k - x_k$ satisfying (3.7) such that

$$(3.9) \quad \|s\|_k \geq \min \left\{ \Delta_k, \kappa_d \frac{\|\nabla m_k(x_k)\|_k}{\kappa_H} \right\},$$

and hence the size of this step is bounded from zero if $\|\nabla m_k(x_k)\|_k$ and Δ_k are.

Reluctance to use nonpolynomial models in practice can be attributed to the difficulty of solving the subproblem (3.2). However, using the sufficient decrease guaranteed by Lemma 3.1, we are still able to guarantee convergence to first-order critical points. This result is independent of the number of local or global minima that the subproblem may have because of using multimodal models.

Further, we assume that the twice continuously differentiable model used in practice will have first- and second-order derivatives available to solve (3.2). Using a more sophisticated solver may be especially attractive when this expense is negligible relative to evaluation of f at the subproblem solution.

We now state the convergence result for our models of interest and Algorithm 1.

THEOREM 3.2. Suppose that the following two assumptions hold:

(AF) $f \in C^1[\Omega]$ for some open $\Omega \supset \mathcal{L}(x_0)$ (with $\mathcal{L}(x_0)$ defined in (3.3)), ∇f is Lipschitz continuous on $\mathcal{L}(x_0)$, and f is bounded on $\mathcal{L}(x_0)$.

(AM) For all $k \geq 0$ we have $m_k \in C^2[\mathcal{B}_k]$, $\infty > \kappa_H \geq \max_{x \in \mathcal{B}_k} \|\nabla^2 m_k(x)\|$, and m_k can be made (and verified to be) fully linear by some finite procedure.

Then for the sequence of iterates generated by Algorithm 1, we have

$$(3.10) \quad \lim_{k \rightarrow \infty} \nabla f(x_k) = 0.$$

Proof. This follows in large part from the lemmas in [9] with minor changes made to accommodate our sufficient decrease condition and the trust-region norm employed. These lemmas, and further explanation where needed, are provided in [36]. \square

4. Radial Basis Functions and ORBIT. Having outlined the fundamental conditions in Theorem 3.2 needed to show convergence of Algorithm 1, in this section we analyze which radial basis function models satisfy these conditions. We also show how the ORBIT algorithm fits in this globally convergent framework.

Throughout this section we drop the dependence of the model on the iteration number, but we intend for the model m and base point x_b to be the k th model and iterate, m_k and x_k , in the trust-region algorithm of the previous section.

An alternative to polynomials is an interpolating surrogate that is a linear combination of nonlinear nonpolynomial basis functions. One such model is of the form

$$(4.1) \quad m(x_b + s) = \sum_{j=1}^{|\mathcal{Y}|} \lambda_j \phi(\|s - y_j\|) + P(s),$$

where $\phi : \mathbb{R}_+ \rightarrow \mathbb{R}$ is a univariate function and $P \in \mathcal{P}_{d-1}^n$ is a polynomial as in section 2. Such models are called *radial basis functions* because $m(x_b + s) - P(s)$ is a linear combination of shifts of a function that is constant on spheres in \mathbb{R}^n .

Interpolation by RBFs on scattered data has only recently gained popularity in practice [5]. In the context of optimization, RBF models have been used primarily for global optimization [4, 16, 33] because they are able to model multimodal/nonconvex functions and interpolate a large number of points in a numerically stable manner.

To our knowledge, Oeuvray was the first to employ RBFs in a local optimization algorithm. In his 2005 dissertation [29], he introduced **BOOSTERS**, a derivative-free trust-region algorithm using a cubic RBF model with a linear tail. Oeuvray was motivated by medical image registration problems and was particularly interested in “doping” his algorithm with gradient information [30]. When the number of interpolation points is fixed from one iteration to the next, Oeuvray also showed that the RBF model parameters λ and ν can be updated in the same complexity as the underdetermined quadratics from [32] (interpolating the same number of points).

4.1. Conditionally Positive Definite Functions. We now define the fundamental property we rely on, using the notation of Wendland [35].

DEFINITION 4.1. Let π be a basis for \mathcal{P}_{d-1}^n , with the convention that $\pi = \emptyset$ if $d = 0$. A function ϕ is said to be conditionally positive definite of order d if for all sets of distinct points $\mathcal{Y} \subset \mathbb{R}^n$ and all $\lambda \neq 0$ satisfying $\sum_{j=1}^{|\mathcal{Y}|} \lambda_j \pi(y_j) = 0$, the quadratic form $\sum_{i,j=1}^{|\mathcal{Y}|} \lambda_i \lambda_j \phi(\|y_i - y_j\|)$ is positive.

Table 4.1 lists examples of popular radial functions and their orders of conditional positive definiteness. Note that if a radial function ϕ is conditionally positive definite of order d , then it is also conditionally positive definite of order $\hat{d} \geq d$ [35, p. 98].

TABLE 4.1

Popular twice continuously differentiable RBFs and order of conditional positive definiteness.

$\phi(r)$	Order	Parameters	Example
r^β	2	$\beta \in (2, 4)$	Cubic, r^3
$(\gamma^2 + r^2)^\beta$	2	$\gamma > 0, \beta \in (1, 2)$	Multiquadric I, $(\gamma^2 + r^2)^{3/2}$
$-(\gamma^2 + r^2)^\beta$	1	$\gamma > 0, \beta \in (0, 1)$	Multiquadric II, $-\sqrt{\gamma^2 + r^2}$
$(\gamma^2 + r^2)^{-\beta}$	0	$\gamma > 0, \beta > 0$	Inverse Multiquadric, $(\gamma^2 + r^2)^{-1/2}$
$\exp(-r^2/\gamma^2)$	0	$\gamma > 0$	Gaussian, $\exp(-r^2/\gamma^2)$

We now use the property of conditional positive definiteness to uniquely determine an RBF model that interpolates data on a set \mathcal{Y} . Let $\Phi_{i,j} = \phi(\|y_i - y_j\|)$ define the square matrix $\Phi \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{Y}|}$, and let Π be the polynomial matrix $\Pi_{i,j} = \pi_i(y_j)$ as in section 2 so that $P(s) = \sum_{i=1}^{\hat{p}} \nu_i \pi_i(s)$. Provided that \mathcal{Y} is \mathcal{P}_{d-1}^n -unisolvent (as in Definition 2.1), we have the equivalent nonsingular symmetric linear system:

$$(4.2) \quad \begin{bmatrix} \Phi & \Pi^T \\ \Pi & 0 \end{bmatrix} \begin{bmatrix} \lambda \\ \nu \end{bmatrix} = \begin{bmatrix} f \\ 0 \end{bmatrix}.$$

The top set of equations corresponds to the interpolation conditions in (2.5) for the RBF model in (4.1), while the lower set ensures uniqueness of the solution.

As in section 2 for polynomial models, for conditionally positive definite functions of order d , a sufficient condition for the nonsingularity of (4.2) is that the points in \mathcal{Y} be distinct and yield a Π^T of full column rank. Clearly this condition is *geometric*, depending only on the location of (but not function values at) the data points.

The saddle point problem in (4.2) will generally be indefinite. However, we employ a null-space method that directly relies on the conditional positive definiteness of ϕ . If Π^T is full rank, then $R \in \mathbb{R}^{(n+1) \times (n+1)}$ is nonsingular from the truncated QR factorization $\Pi^T = QR$. By the lower set of equations in (4.2) we must have $\lambda = Z\omega$ for $\omega \in \mathbb{R}^{|\mathcal{Y}| - n - 1}$ and any orthogonal basis Z for $\mathcal{N}(\Pi)$. Hence (4.2) reduces to

$$(4.3) \quad Z^T \Phi Z \omega = Z^T f,$$

$$(4.4) \quad R\nu = Q^T(f - \Phi Z\omega).$$

Given that Π^T is full rank and the points in \mathcal{Y} are distinct, Definition 4.1 directly implies that $Z^T \Phi Z$ is positive definite for any ϕ that is conditionally positive definite of at most order d . Positive definiteness of $Z^T \Phi Z$ guarantees the existence of a nonsingular lower triangular Cholesky factor L such that

$$(4.5) \quad Z^T \Phi Z = LL^T,$$

and the isometry of Z gives the bound

$$(4.6) \quad \|\lambda\| = \|ZL^{-T}L^{-1}Z^T f\| \leq \|L^{-1}\|^2 \|f\|.$$

4.2. Fully Linear RBF Models. Thus far we have maintained a very general RBF framework. For the convergence results in section 3 to apply, we now focus on a more specific set of radial functions that satisfy two additional conditions:

- $\phi \in C^2[\mathbb{R}_+]$ and $\phi'(0) = 0$.
- ϕ is conditionally positive definite of order 2 or less.

The first condition ensures that the resulting RBF model is twice continuously differentiable. The second condition is useful for restricting ourselves to models of the form (4.1) with a linear tail $P \in \mathcal{P}_1^n$.

For RBF models that are twice continuously differentiable and have a linear tail,

$$(4.7) \quad \nabla m(x_b + s) = \sum_{\{y_i \in \mathcal{Y}: y_i \neq s\}} \lambda_i \phi'(\|s - y_i\|) \frac{s - y_i}{\|s - y_i\|} + \nabla P(s),$$

$$(4.8) \quad \nabla^2 m(x_b + s) = \sum_{y_i \in \mathcal{Y}} \lambda_i \Theta(\|s - y_i\|),$$

with

$$(4.9) \quad \Theta(r) = \begin{cases} \frac{\phi'(\|r\|)}{\|r\|} I_n + \left(\phi''(\|r\|) - \frac{\phi'(\|r\|)}{\|r\|} \right) \frac{r}{\|r\|} \frac{r^T}{\|r\|} & \text{if } r \neq 0, \\ \phi''(0) I_n & \text{if } r = 0, \end{cases}$$

where we have explicitly defined these derivatives for the special case when s is one of the interpolation knots in \mathcal{Y} .

The following lemma is a consequence of an unproven statement in Oeuvray's dissertation [29], which we could not otherwise locate in the literature. It provides necessary and sufficient conditions on ϕ for the RBF model m to be twice continuously differentiable.

LEMMA 4.2. *The model m defined in (4.1) is twice continuously differentiable on \mathbb{R}^n if and only if $\phi \in C^2[\mathbb{R}_+]$ and $\phi'(0) = 0$.*

Proof. We begin by noting that the polynomial tail P and composition with the sum over \mathcal{Y} are both smooth. Moreover, away from any of the points in \mathcal{Y} , m is clearly twice continuously differentiable if and only if $\phi \in C^2[\mathbb{R}_+]$. It now remains only to treat the case when $s = y_j \in \mathcal{Y}$.

If ϕ' is continuous but $\phi'(0) \neq 0$, then since $\frac{s - y_j}{\|s - y_j\|}$ is always of bounded magnitude but does not exist as $s \rightarrow y_j$, we have that ∇m in (4.7) is not continuous at y_j .

We conclude by noting that $\phi'(0) = 0$ is sufficient for the continuity of $\nabla^2 m$ at y_j . To see this, recall from L'Hôpital's rule in calculus that $\lim_{a \rightarrow 0} \frac{g(a)}{a} = g'(0)$, provided $g(0) = 0$ and g is differentiable at 0. Applying this result with $g = \phi'$, we have that

$$\lim_{s \rightarrow y_j} \frac{\phi'(\|s - y_j\|)}{\|s - y_j\|} = \phi''(0).$$

Hence the second term in the expression for Θ in (4.9) vanishes as $r \rightarrow 0$, leaving only the first; that is, $\lim_{r \rightarrow 0} \Theta(r) = \phi''(0) I_n$ exists. \square

We note that this result implies that models using the thin-plate spline radial function $\phi(r) = r^2 \log(r)$ are not twice continuously differentiable and hence do not fit in our framework.

Having established conditions for the twice differentiability of the radial portion of m in (4.1), we now focus on the linear tail P . Without loss of generality, we assume that the base point x_b is an interpolation point so that $y_1 = 0 \in \mathcal{Y}$. Employing the standard linear basis and permuting the points, we then have that the polynomial matrix $\Pi_{i,j} = \pi_i(y_j)$ is of the form

$$(4.10) \quad \Pi = \begin{bmatrix} Y & 0 & y_{n+2} & \cdots & y_{|\mathcal{Y}|} \\ e^T & 1 & 1 & \cdots & 1 \end{bmatrix},$$

TABLE 4.2

Upper bounds on RBF components (assumes $\gamma > 0$, $r \in [0, \Delta]$, β as in Table 4.1).

$\phi(r)$	$ \phi(r) $	$\left \frac{\phi'(r)}{r}\right $	$ \phi''(r) $
r^β	Δ^β	$\beta\Delta^{\beta-2}$	$\beta(\beta-1)\Delta^{\beta-2}$
$(\gamma^2 + r^2)^\beta$	$(\gamma^2 + \Delta^2)^\beta$	$2\beta(\gamma^2 + \Delta^2)^{\beta-1}$	$2\beta(\gamma^2 + \Delta^2)^{\beta-1} \left(1 + \frac{2(\beta-1)\Delta^2}{\gamma^2 + \Delta^2}\right)$
$-(\gamma^2 + r^2)^\beta$	$(\gamma^2 + \Delta^2)^\beta$	$2\beta\gamma^{2(\beta-1)}$	$2\beta\gamma^{2(\beta-1)}$
$(\gamma^2 + r^2)^{-\beta}$	$\gamma^{-2\beta}$	$2\beta\gamma^{-2(\beta+1)}$	$2\beta\gamma^{-2(\beta+1)}$
$\exp(-r^2/\gamma^2)$	1	$2/\gamma^2$	$2/\gamma^2$

where e is the vector of ones and Y denotes a matrix of n particular nonzero points in \mathcal{Y} .

Recall that, in addition to the distinctness of the points in \mathcal{Y} , a condition for the nonsingularity of the RBF system (4.2) is that the first $n+1$ columns of Π in (4.10) are linearly independent. This is exactly the condition needed for the fully linear interpolation models in section 2, where bounds for the matrix Y were provided.

To fit RBF models with linear tails into the globally convergent trust-region framework of section 3, we need only to show that the model Hessians are bounded by some fixed constant κ_H .

From (4.8) and (4.9), the magnitude of the Hessian depends only on the quantities λ , $\left|\frac{\phi'(r)}{r}\right|$, and $|\phi''(r)|$. As an example, Table 4.2 provides bounds on the last two quantities for the radial functions in Table 4.1 when r is restricted to lie in the interval $[0, \Delta]$. In particular, these bounds provide an upper bound for

$$(4.11) \quad h_\phi(\Delta) = \max \left\{ 2 \left| \frac{\phi'(r)}{r} \right| + |\phi''(r)| : r \in [0, \Delta] \right\}.$$

From (4.6) we also have a bound on λ provided that the appropriate Cholesky factor L is of bounded norm. We bound $\|L^{-1}\|$ inductively by building up the interpolation set \mathcal{Y} one point at a time. This inductive method lends itself well to a practical implementation and was inspired by the development in [4].

To start this inductive argument, we assume that \mathcal{Y} consists of $n+1$ points that are \mathcal{P}_1^n -unisolvent. With only these $n+1$ points, $\lambda = 0$ is the unique solution to (4.2), and hence the RBF model is linear. To include an additional point $y \in \mathbb{R}^n$ in the interpolation set \mathcal{Y} (beyond the initial $n+1$ points), we appeal to the following lemma (derived in [38]).

LEMMA 4.3. *Let \mathcal{Y} be such that Π is full rank and $LL^T = Z^T\Phi Z$ is invertible as in (4.5). If $y \in \mathbb{R}^n$ is added to \mathcal{Y} , then the new Cholesky factor L_y has an inverse*

$$(4.12) \quad L_y^{-1} = \begin{bmatrix} L^{-1} & 0 \\ \frac{-v_y^T L^{-T} L^{-1}}{\tau(y)} & \frac{1}{\tau(y)} \end{bmatrix}, \quad \text{with } \tau(y) = \sqrt{\sigma_y - \|L^{-1}v_y\|^2},$$

provided that the constant $\tau(y)$ is positive.

Here we see that only the last row of L_y^{-1} is affected by the addition of the new point y . As noted in [38], the constant σ_y and vector v_y in Lemma 4.3 appear in the reduced $Z_y^T \Phi_y Z_y = L_y L_y^T$ when y is added, and can be obtained by applying $n+1$ Givens rotations to Π_y^T . The following lemma bounds the resulting Cholesky factor L_y^{-1} as a function of the previous factor L^{-1} , v_y , and $\tau(y)$.

Algorithm 2 Algorithm for adding additional interpolation points.

2.0. Input $\mathcal{D} = \{d_1, \dots, d_{|\mathcal{D}|}\} \subset \mathbb{R}^n$, \mathcal{Y} consisting of $n+1$ sufficiently affinely independent points, constants $\theta_2 > 0$, $\Delta > 0$, and $p_{\max} \geq n+1$.

2.1. Using \mathcal{Y} , compute the Cholesky factorization $LL^T = Z^T \Phi Z$ as in (4.5).

2.2. For all $y \in \mathcal{D}$ such that $\|y\|_k \leq \Delta$:

If $\tau(y) \geq \theta_2$,

Update $\mathcal{Y} \leftarrow \mathcal{Y} \cup \{y\}$, $Z \leftarrow Z_y$, $L \leftarrow L_y$,

If $|\mathcal{Y}| = p_{\max}$, **return**.

LEMMA 4.4. *If $\|L^{-1}\| \leq \kappa$ and $\tau(y) \geq \theta > 0$, then*

$$(4.13) \quad \|L_y^{-1}\|^2 \leq \kappa + \frac{1}{\theta^2} (1 + \|v_y\| \kappa^2)^2.$$

Proof. Let $w_y = (w, \tilde{w}) \in R^{|\mathcal{Y}|+1}$ be an arbitrary vector with $\|w_y\| = 1$. Then

$$\begin{aligned} \|L_y^{-1} w_y\|^2 &= \|L^{-1} w\|^2 + \frac{1}{\tau(y)^2} (\tilde{w} - v_y^T L^{-T} L^{-1} w)^2 \\ &\leq \kappa + \frac{1}{\theta^2} \left(\tilde{w}^2 - 2\tilde{w} v_y^T L^{-T} L^{-1} w + (v_y^T L^{-T} L^{-1} w)^2 \right) \\ &\leq \kappa + \frac{1}{\theta^2} \left(1 + 2 \|L^{-1} v_y\| \|L^{-1} w\| + (\|L^{-1} v_y\| \|L^{-1} w\|)^2 \right) \\ &\leq \kappa + \frac{1}{\theta^2} (1 + \|v_y\| \kappa^2)^2. \quad \square \end{aligned}$$

Lemma 4.4 suggests the procedure given in Algorithm 2, which we use in **ORBIT** to iteratively add previously evaluated points to the interpolation set \mathcal{Y} . Before this algorithm is called, we assume that \mathcal{Y} consists of $n+1$ sufficiently affinely independent points generated as described in section 2 and hence the initial L matrix is empty.

Figure 4.1 (a) gives an example of $\tau(y)^{-1}$ values for different interpolation sets in \mathbb{R}^2 . In particular we note that $\tau(y)$ approaches zero as y approaches any of the points already in the interpolation set \mathcal{Y} . Figure 4.1 (b) shows the behavior of $\|L_y^{-1}\|^2$ for the same interpolation sets and illustrates the relative correspondence between the values of $\tau(y)^{-1}$ and $\|L_y^{-1}\|^2$.

We now assume that both \mathcal{Y} and the point y being added to the interpolation set belong to some bounded domain $\{x \in \mathbb{R}^n : \|x\|_k \leq \Delta\}$. Thus the quantities $\{\|x - z\| : x, z \in \mathcal{Y} \cup y\}$ are all of magnitude no more than $2c_1\Delta$, since $\|\cdot\| \leq c_1 \|\cdot\|_k$. The elements in $\Phi_{i,j} = \phi(\|y_i - y_j\|)$ and $\phi_y = [\phi(\|y - y_1\|), \dots, \phi(\|y - y_{|\mathcal{Y}|}\|)]^T$ are bounded by $k_\phi(2c_1\Delta)$, where

$$(4.14) \quad k_\phi(2c_1\Delta) = \max\{|\phi(r)| : r \in [0, 2c_1\Delta]\}.$$

Bounds for the specific ϕ functions of the radial basis functions of interest are provided in Table 4.2. Using the isometry of Z_y we hence have the bound

$$(4.15) \quad \|v_y\| \leq \sqrt{|\mathcal{Y}|(|\mathcal{Y}| + 1)} k_\phi(2c_1\Delta),$$

independent of where in $\{x \in \mathbb{R}^n : \|x\|_k \leq \Delta\}$ the point y lies, which can be used in (4.13) to bound $\|L_y^{-1}\|$. The following theorem gives the resulting bound.

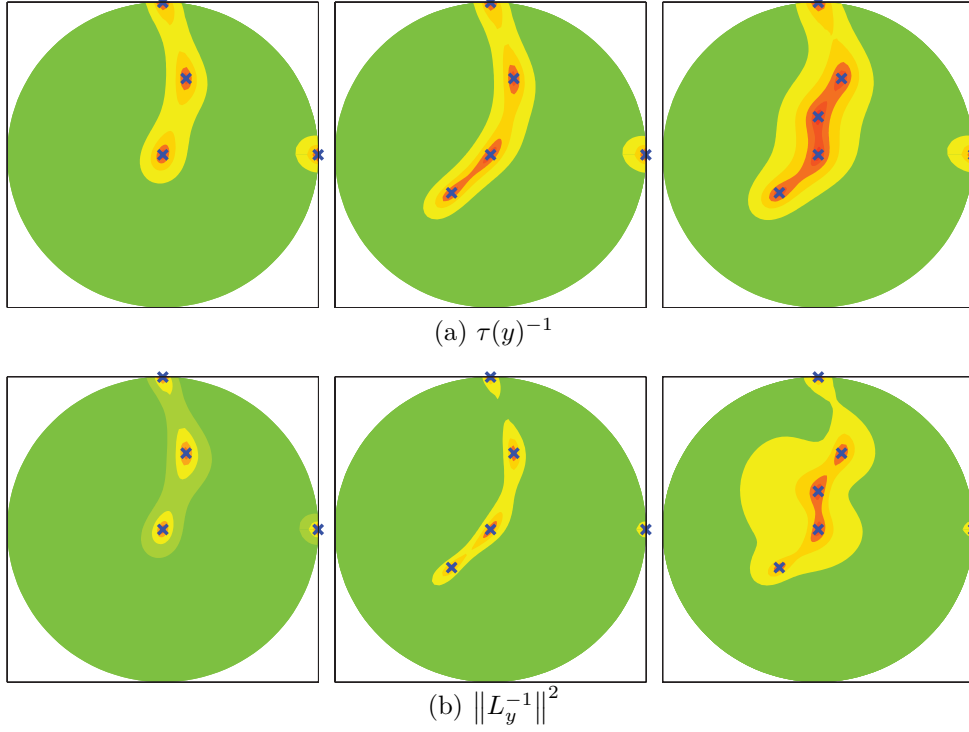


FIG. 4.1. Contours for $\tau(y)^{-1}$ and $\|L_y^{-1}\|^2$ values (4.12) for a multiquadric RBF interpolating 4, 5, and 6 points in \mathbb{R}^2 (log scale). The quantities grow as the interpolation points are approached.

THEOREM 4.5. Let $\mathcal{B} = \{x \in \mathbb{R}^n : \|x - x_b\|_k \leq \Delta\}$. Let $\mathcal{Y} \subset \mathcal{B} - x_b$ be a set of distinct interpolation points, $n + 1$ of which are affinely independent and $|f(x_b + y_i)| \leq f_{\max}$ for all $y_i \in \mathcal{Y}$. Then for a model of the form (4.1), with a bound h_ϕ as defined in (4.11), interpolating f on $x_b + \mathcal{Y}$, we have that for all $x \in \mathcal{B}$

$$(4.16) \quad \|\nabla^2 m(x)\| \leq |\mathcal{Y}| \|L^{-1}\|^2 h_\phi(2c_1\Delta) f_{\max} =: \kappa_H.$$

Proof. Let $r_i = s - y_i$, and note that when s and \mathcal{Y} both belong to $\mathcal{B} - x_b$, $\|r_i\| \leq c_1 \|r_i\|_k \leq 2c_1\Delta$ for $i = 1, \dots, |\mathcal{Y}|$. Thus for an arbitrary w with $\|w\| = 1$,

$$\begin{aligned} \|\nabla^2 m(x_b + s)w\| &\leq \sum_{i=1}^{|\mathcal{Y}|} |\lambda_i| \left\| \frac{\phi'(\|r_i\|)}{\|r_i\|} w + \left(\phi''(\|r_i\|) - \frac{\phi'(\|r_i\|)}{\|r_i\|} \right) \frac{r_i^T w}{\|r_i\|} \frac{r_i}{\|r_i\|} \right\|, \\ &\leq \sum_{i=1}^{|\mathcal{Y}|} |\lambda_i| \left[2 \left| \frac{\phi'(\|r_i\|)}{\|r_i\|} \right| + |\phi''(\|r_i\|)| \right] \\ &\leq \|\lambda\|_1 h(2c_1\Delta) \leq \sqrt{|\mathcal{Y}|} \|L^{-1}\|^2 \|f\| h(2c_1\Delta), \end{aligned}$$

where the last two inequalities follow from (4.11) and (4.6), respectively. Noting that $\|f\| \leq \sqrt{|\mathcal{Y}|} f_{\max}$ gives the desired result. \square

4.3. RBF Models in ORBIT. Having shown how RBFs fit into the globally convergent framework for fully linear models, we collect some final details of ORBIT, consisting of Algorithm 1 and the RBF model formation summarized in Algorithm 3.

Algorithm 3 Algorithm for constructing model m_k .

- 3.0.** Input $\mathcal{D} \subset \mathbb{R}^n$, constants $\theta_2 > 0$, $\theta_4 \geq \theta_3 \geq 1$, $\theta_1 \in (0, \frac{1}{\theta_3}]$, $\Delta_{\max} \geq \Delta_k > 0$, and $p_{\max} \geq n + 1$.
- 3.1.** Seek affinely independent interpolation set \mathcal{Y} within distance $\theta_3 \Delta_k$.
 Save z_1 as a model-improving direction for use in Step 1.5 of Algorithm 1.
 If $|\mathcal{Y}| < n + 1$ (and hence m_k is not fully linear):
 Seek $n + 1 - |\mathcal{Y}|$ additional points in \mathcal{Y} within distance $\theta_4 \Delta_{\max}$.
 If $|\mathcal{Y}| < n + 1$, evaluate f at remaining $n + 1 - |\mathcal{Y}|$ model points so that $|\mathcal{Y}| = n + 1$.
- 3.2.** Use up to $p_{\max} - n - 1$ additional points in \mathcal{D} within $\theta_4 \Delta_{\max}$ using Algorithm 2.
- 3.3.** Obtain model parameters by (4.3) and (4.4).

Algorithm 3 requires that the interpolation points in \mathcal{Y} lie within some constant factor of the largest trust region Δ_{\max} . This region, $\mathcal{B}_{\max} = \{y \in \mathbb{R}^n : \|y\|_k \leq \theta_4 \Delta_{\max}\}$, is chosen to be larger than the current trust region so that the algorithm can make use of more points previously evaluated in the course of the optimization.

In Algorithm 3 we certify a model to be fully linear if $n + 1$ points within $\{y \in \mathbb{R}^n : \|y\|_k \leq \theta_3 \Delta_k\}$ result in pivots larger than θ_1 , where the constant θ_1 is chosen so as to be attainable by the model directions (scaled by Δ_k) discussed in section 2.

If not enough points are found, the model will not be fully linear; thus, we must expand the search for affinely independent points within the larger region \mathcal{B}_{\max} . If still fewer than $n + 1$ points are available, we must evaluate f along a set of the model-improving directions Z to ensure that \mathcal{Y} is \mathcal{P}_1^n -unisolvant.

Additional available points within \mathcal{B}_{\max} are added to the interpolation set \mathcal{Y} provided that they keep $\tau(y) \geq \theta_2 > 0$, until a maximum of p_{\max} points are in \mathcal{Y} .

Since we have assumed that f is bounded on $\mathcal{L}(x_0)$ and that $\mathcal{Y} \subset \mathcal{B}_{\max}$, the bound (4.16) holds for all models used by the algorithm, regardless of whether they are fully linear. Provided that the radial function ϕ is chosen to satisfy the requirements of Lemma 4.2, m will be twice continuously differentiable. Hence ∇m is Lipschitz continuous on \mathcal{B}_{\max} , and κ_H in (3.8) is one possible Lipschitz constant. When combined with the results of section 2 showing that such interpolation models can be made fully linear in a finite procedure, Theorem 3.2 guarantees that $\lim_{k \rightarrow \infty} \nabla f(x_k) = 0$ for trust-region algorithms using these RBFs, and ORBIT in particular.

5. Computational Experiments. We now present numerical results aimed at determining the effect of selecting different types of RBF models. We follow the benchmarking procedures in [26], with the derivative-free convergence test

$$(5.1) \quad f(x_0) - f(x) \geq (1 - \tau)(f(x_0) - f_L),$$

where $\tau > 0$ is a tolerance, x_0 is the starting point, and f_L is the smallest value of f obtained by any tested solver within a fixed number, μ_f , of function evaluations. We note that in (5.1), a problem is “solved” when the achieved reduction from the initial value, $f(x_0) - f(x)$, is at least $1 - \tau$ times the best possible reduction, $f(x_0) - f_L$.

For each solver $s \in \mathcal{S}$ and problem $p \in \mathcal{P}$, we define $t_{p,s}$ as the number of function evaluations required by s to satisfy the convergence test (5.1) on p , with the convention that $t_{p,s} = \infty$ if s does not satisfy the convergence test on p within μ_f evaluations.

If we assume that (i) the differences in times for solvers to determine a point for evaluation of $f(x)$ are negligible relative to the time to evaluate the function and (ii)

the function requires the same amount of time to evaluate at any point in its domain, then differences in the measure $t_{p,s}$ roughly correspond to differences in computing time. Assumption (i) is reasonable for the computationally expensive simulation-based problems motivating this work.

Given this measure, we define the *data profile* $d_s(\alpha)$ for solver $s \in \mathcal{S}$ as

$$(5.2) \quad d_s(\alpha) = \frac{1}{|\mathcal{P}|} \left| \left\{ p \in \mathcal{P} : \frac{t_{p,s}}{n_p + 1} \leq \alpha \right\} \right|,$$

where n_p is the number of variables in problem $p \in \mathcal{P}$. We note that the data profile $d_s : \mathbb{R} \rightarrow [0, 1]$ is a nondecreasing step function independent of the data profiles of the other solvers $\mathcal{S} \setminus \{s\}$, provided that f_L is fixed. By this definition, $d_s(\kappa)$ is the fraction of problems that can be solved within κ *simplex gradient estimates* (and hence a budget of $\kappa(n_p + 1)$ function evaluations).

5.1. Smooth Test Problems. We begin by considering the test set \mathcal{P}_S of 53 smooth, nonlinear least-squares problems in [26]. Each unconstrained problem is defined by a starting point x_0 and a function $f(x) = \sum_{i=1}^k f_i(x)^2$, comprising a set of smooth components. The functions vary in dimension from $n = 2$ to $n = 12$, with the 53 problems being distributed roughly uniformly across these dimensions. The maximum number of function evaluations is set to $\mu_f = 1300$ so that at least the equivalent of 100 simplex gradient estimates can be obtained on all the problems in \mathcal{P}_S . The initial trust-region radius is set to $\Delta_0 = \max\{1, \|x_0\|_\infty\}$ for each problem.

The ORBIT implementation illustrated here relies on a 2-norm trust region with parameter values as in [38]: $\eta_0 = 0$, $\eta_1 = .2$, $\gamma_0 = \frac{1}{2}$, $\gamma_1 = 2$, $\Delta_{\max} = 10^3 \Delta_0$, $\epsilon = 10^{-10}$, $\kappa_d = 10^{-4}$, $\alpha = .9$, $\mu = 2000$, $\beta = 1000$, $\theta_1 = 10^{-3}$, $\theta_2 = 10^{-7}$, $\theta_3 = 10$, $\theta_4 = \max(\sqrt{n}, 10)$. In addition to the backtracking line search detailed here, we use an augmented Lagrangian method to approximately solve the trust-region subproblem.

The first solver set we consider is the set \mathcal{S}_A consisting of four radial basis function types for ORBIT:

Multiquadric: $\phi(r) = -\sqrt{1 + r^2}$, with $p_{\max} = 2n + 1$.

Cubic: $\phi(r) = r^3$, with $p_{\max} = 2n + 1$.

Gaussian: $\phi(r) = \exp(-r^2)$, with $p_{\max} = 2n + 1$.

Thin plate: $\phi(r) = r^2 \log(r)$, with $p_{\max} = 2n + 1$.

The common theme among these models is that they interpolate at most $p_{\max} = 2n + 1$ points, chosen because this is the number of interpolation points recommended by Powell for the NEWUOA algorithm [32]. We tested other values of the parameter γ used by multiquadric and Gaussian RBFs but found that $\gamma = 1$ worked well for both.

In our testing, we examined accuracy levels of $\tau = 10^{-k}$ for several k . For the sake of brevity, in Figure 5.1 we present the data profiles for $k = 1$ and $k = 5$. Recall that $\tau = 0.1$ corresponds to a 90% reduction relative to the best possible reduction in $\mu_f = 1300$ function evaluations. As discussed in [26], data profiles are used to see which solver is likely to achieve a given reduction of the function within a specific computational budget. For example, given the equivalent of 15 simplex gradients ($15(n + 1)$ function evaluations), we see that the cubic, multiquadric, Gaussian, and thin-plate spline variants respectively solve 38%, 30%, 27%, and 30% of problems to $\tau = 10^{-5}$ accuracy.

For the accuracy levels shown, the cubic variant is generally best (especially given small budgets), while the Gaussian and thin-plate spline variants are generally worst. The differences are smaller than those seen in [26], where \mathcal{S} consisted of three very different solvers.

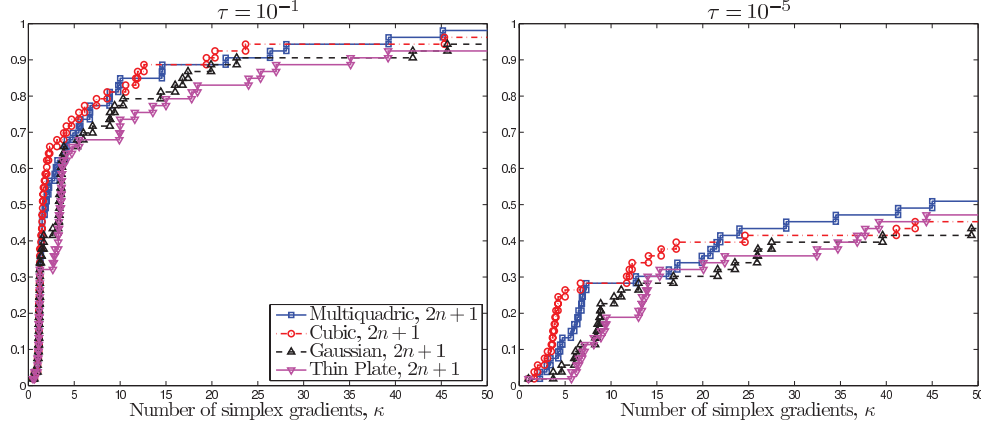


FIG. 5.1. Data profiles $d_s(\kappa)$ for different RBF types with $p_{\max} = 2n+1$ on the smooth problems \mathcal{P}_S . These profiles show the fraction of problems solved as a function of a computational budget of simplex gradients ($\kappa(n+1)$ function evaluations).

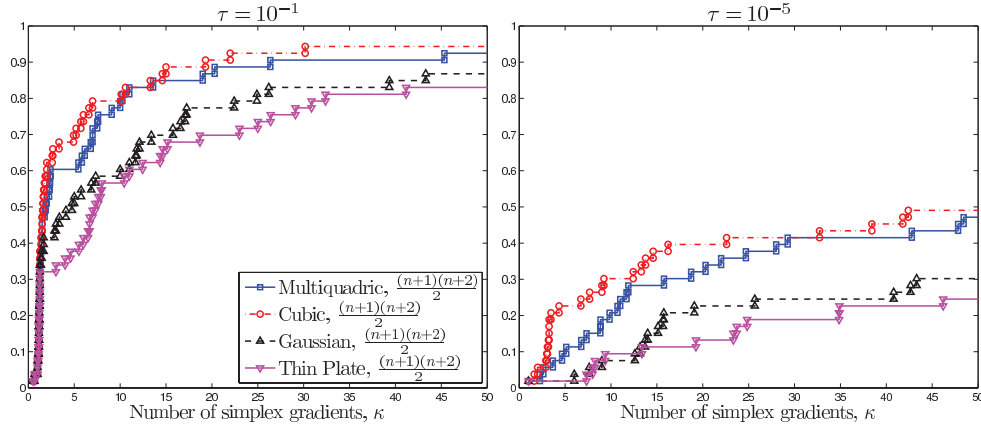


FIG. 5.2. Data profiles $d_s(\kappa)$ for different RBF types with $p_{\max} = \frac{(n+1)(n+2)}{2}$ on the smooth problems \mathcal{P}_S . These profiles show the fraction of problems solved as a function of a computational budget of simplex gradients ($\kappa(n+1)$ function evaluations).

The second solver set, \mathcal{S}_B , consists of the same four radial basis function types:

Multiquadric: $\phi(r) = -\sqrt{1+r^2}$, with $p_{\max} = \frac{(n+1)(n+2)}{2}$.

Cubic: $\phi(r) = r^3$, with $p_{\max} = \frac{(n+1)(n+2)}{2}$.

Gaussian: $\phi(r) = \exp(-r^2)$, with $p_{\max} = \frac{(n+1)(n+2)}{2}$.

Thin plate: $\phi(r) = r^2 \log(r)$, with $p_{\max} = \frac{(n+1)(n+2)}{2}$.

Here, the maximum number of points being interpolated corresponds to the number of points needed to uniquely fit an interpolating quadratic model. This choice is made solely to indicate how the performance changes with a larger number of interpolation points.

Figure 5.2 shows the data profiles for the accuracy levels $\tau \in \{10^{-1}, 10^{-5}\}$. The cubic variant is again generally best (especially given small budgets), but there are now larger differences among the variants. When the equivalent of 15 simplex gradients is available, we see that the cubic, multiquadric, Gaussian, and thin-plate spline

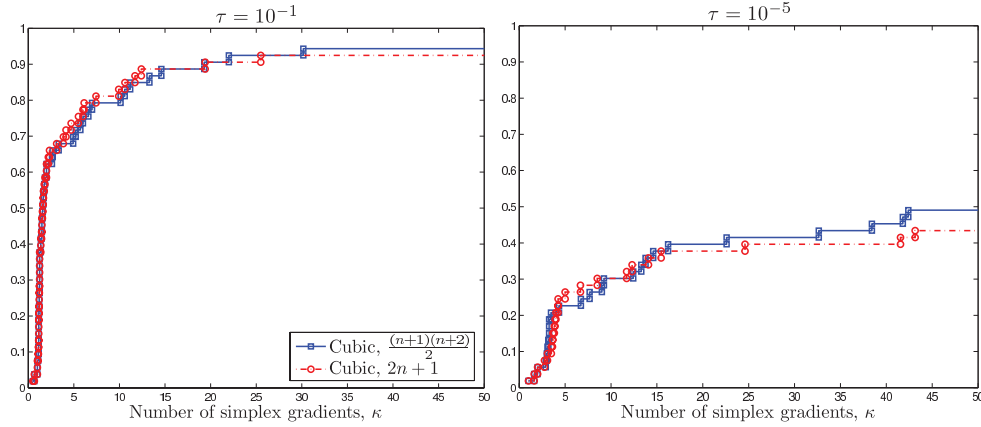


FIG. 5.3. Effect of changing the maximum number of interpolation points, p_{\max} , on the data profiles $d_s(\kappa)$ for the smooth problems \mathcal{P}_S .

variants are respectively able to now solve 37%, 28%, 16%, 11% of problems to an accuracy level of $\tau = 10^{-5}$. We note that the raw data in Figure 5.2 should not be quantitatively compared with that in Figure 5.1 because the best function value found for each problem is obtained from only the solvers tested (in \mathcal{S}_A or \mathcal{S}_B) and hence the convergence tests differ.

Our final test on these problems compares the best variants for the two different maximum numbers of interpolation points. The solver set \mathcal{S}_C consists of the following:

Cubic A: $\phi(r) = r^3$, with $p_{\max} = 2n + 1$.

Cubic B: $\phi(r) = r^3$, with $p_{\max} = \frac{(n+1)(n+2)}{2}$.

Figure 5.3 shows that these two variants perform comparably, with differences smaller than those seen in Figures 5.1 and 5.2. As expected, as the number of function evaluations grows, the variant that is allowed to interpolate more points performs better. This variant also performs better when higher accuracy levels are demanded, and we attribute this performance to the fact that the model interpolating more points is generally a better approximation of the function f . The main downside of interpolating more points is that the linear systems in section 4 will also grow, resulting in a higher linear algebra cost per iteration. As we will see in the next set of tests, for many applications, this cost may be viewed as negligible relative to the cost of evaluating the function f .

We are, however, surprised to see that the $2n + 1$ variant performs better for some smaller budgets. For example, this variant performs slightly better between 5 and 15 simplex gradient estimates when $\tau = 10^{-1}$ and between 4 and 9 simplex gradient estimates when $\tau = 10^{-5}$. Since the initial $n + 1$ evaluations are common to both variants and since the parameter p_{\max} has no effect on the subroutine determining the sufficiently affinely independent points, we might expect that the variant interpolating more points would do at least as well as the variant interpolating fewer points.

Further results comparing ORBIT (in 2-norm and ∞ -norm trust regions) with NEWUOA on a set of noisy test problems are provided in [38].

5.2. An Environmental Application. We now illustrate the use of RBF models on a computationally expensive application problem.

The Blaine Naval Ammunition Depot is a 48,800 acre site east of Hastings, Ne-

braska. Nearly half of the naval ammunition used in World War II was produced here, with the result that much toxic waste was generated and disposed of on the site. Both trichloroethylene (TCE), a probable carcinogen, and trinitrotoluene (TNT), a possible carcinogen, are present in the groundwater. Because of their possible connection to cancer, the two chemicals have EPA allowable limits that are extremely low. For the problem considered here the remaining concentrations (after 30 years of remediation is completed) for TCE must be less than 5.0 parts per billion (ppb) and for TNT must be less than 2.5 ppb. The transport of these contaminants is made especially difficult because the predominant groundwater flow direction changes during the region's irrigation season, which lasts roughly two and a half months.

As part of a collaboration [3, 39] among environmental consultants, academic institutions, and governmental agencies, several optimization problems were formulated. Here we focus on one of the simpler formulations, where we have control over 15 injection and extraction wells located at fixed positions in the site. At each of these wells we can either inject clean water or extract and then treat contaminated water (e.g., with air stripping or carbon adsorption). Each instance of the decision variables hence corresponds to a pumping strategy that will run over a 30-year time horizon. For scaling purposes, each variable is scaled so that the range of realistic pumping rates (corresponding to injection and extraction rates between 0 and 350 gallons per minute) maps to the interval $[0, 1]$.

The groundwater model is formally described in [39] and was calibrated for steady-state and transient conditions to historical plume data based on particle tracking. Groundwater occurs approximately 100 feet below the ground surface. This groundwater is predominantly in a semi-confined layer and is the major water supply source for the region. In the formulation considered, the surface area is partitioned into a model grid of 11,152 cells covering 134 square miles. Six layers (corresponding to the unconfined aquifer, upper confining layer, and the semi-confined aquifer split into four layers) were employed in the vertical direction.

The groundwater flow is simulated using MODFLOW 2000 [34], which solves the transient groundwater equation (see, e.g., [34, p. 10]) using a finite-difference method. Contaminant transport is simulated with the modular three-dimensional transport code MT3D-MS [40]. There are other contaminants present (e.g., forms of PCE, DCE, TCA, and RDX) at the site, but the formulation considered [3, 39] includes these other volatile organic compounds as part of the TCE concentration since they share similar transport behaviors with TCE. Hence, MT3D-MS simulates only the concentration of the two species TCE and TNT.

The objective is to minimize the cost of the pumping strategy (including the electricity needed to run the pumps) plus a financial penalty associated with exceeding the constraints on maximum allowable concentration of TCE and TNT over the 30-year planning horizon. For each pumping strategy, these concentrations are obtained by running the coupled simulators MODFLOW 2000 and MT3D-MS. For a given set of pumping rates, this process required more than 45 minutes on a Pentium 4 dual-core desktop. Thus, each evaluation of the objective function required more than 45 minutes.

In the spirit of [26], in addition to ORBIT we considered three solvers designed to solve unconstrained serial optimization problems using only function values.

NMSMAX is an implementation of the Nelder-Mead method and is due to Higham [17]. We specified that the initial simplex have sides of length Δ_0 . Since NMSMAX is defined for maximization problems, it was given $-f$.

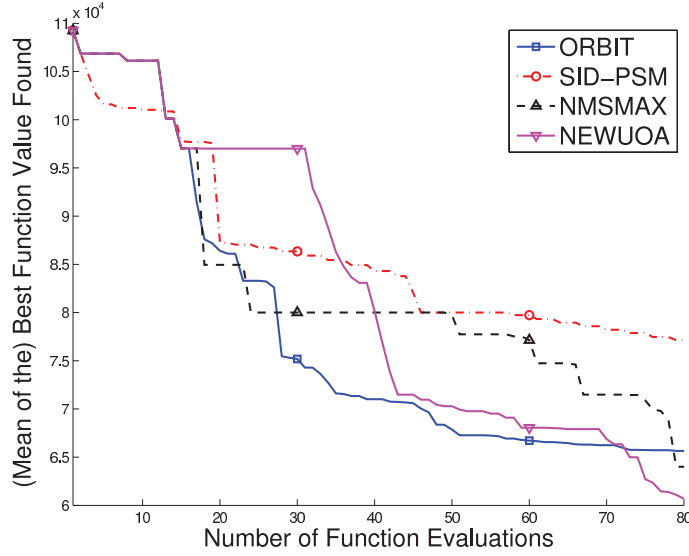


FIG. 5.4. Mean (in 8 trials) of the best function value found for the first 80 evaluations on the Blaine problem. All **ORBIT** runs found a local minimum within 80 evaluations, while **NEWUOA** obtained a lower function value after 72 evaluations.

SID-PSM is a pattern search solver due to Custódio and Vicente [13]. It is especially designed to use previous function evaluations. We used version 0.4 with an initial step size set to Δ_0 . We note that the performance of **SID-PSM** has since been improved with the incorporation of interpolating models (as reported in [12]), but we have reported the originally tested version as an example of an industrial-strength pattern search method not incorporating such models.

NEWUOA is a trust-region solver using a quadratic model and is due to Powell [32]. The number of interpolation points was fixed at the recommended value of $p_{\max} = 2n + 1$, and the initial trust-region radius was set to Δ_0 .

ORBIT used the same parameter values as used on the test functions, with a cubic RBF, initial trust-region radius Δ_0 , and a maximum number of interpolation points taken to be larger than the number of function evaluations, $p_{\max} \geq \mu_f$.

Each of these solvers also requires a starting point x_0 and a maximum number of allowable function evaluations, μ_f . A common selection of $\Delta_0 = 0.1$ was made to standardize the initial evaluations across the collection of solvers. Hence each solver except **SID-PSM** evaluated the same initial $n + 1$ points. **SID-PSM** moves off this initial pattern once it sees a reduction. All other inputs were set to their default values except that we effectively set all termination parameters to zero to ensure that the solvers terminate only after exhausting the budget μ_f function evaluations.

We set $\mu_f = 5(n + 1) = 80$; and since each evaluation (i.e., an environmental model simulation) requires more than 45 minutes, a single run of one solver thus requires nearly 3 CPU-days. As this problem is noisy and has multiple local minima, we chose to run each solver from the same eight starting points generated uniformly at random within the hypercube $[0, 1]^{15}$ of interest. Thus, running four solvers over these eight starting points required roughly 3 CPU-months to obtain.

Figure 5.4 shows the average of the best function value obtained over the course of the first 80 function evaluations. By design, all solvers start from the same function

value. The ORBIT solver does best initially, obtaining a function value of 70,000 in 46 evaluations. The ORBIT trajectory quickly flattens out as it is the first solver to find a local minima, with an average value of 65,600. In this case, however, the local minimum found most quickly by ORBIT has (on average) a higher function value than the point (not yet a local minimum) found by the NEWUOA and NMSMAX solvers after $\mu_f = 80$ evaluations. Hence, in these tests, NEWUOA and NMSMAX are especially good at finding a good minimum for a noisy function. On average, given $\mu_f = 80$ evaluations, NEWUOA finds a point with $f \approx 60,700$. None of these algorithms are designed to be global optimization solvers, so the focus here is on the time to find the first local minimum.

The Blaine problem highlights the fact that solvers will have different performance on different functions and that many application problems contain computational noise and multiple distinct local minima, which can prevent globally convergent local methods from finding good solutions. Comparisons between ORBIT and other derivative-free algorithms on two different problems from environmental engineering can be found in [38]. The results in [38] indicate that two variants of ORBIT outperformed the three other solvers tested on these two environmental problems.

6. Conclusions and Perspectives. In this paper we have introduced and analyzed first-order, derivative-free trust-region algorithms based on radial basis functions, which are globally convergent. We first showed that, provided a function and a model are sufficiently smooth, interpolation on a set of sufficiently affinely independent points is enough to guarantee Taylor-like error bounds for both the model and its gradient. In section 3 we extended the recent derivative-free trust-region framework in [9] to include nonlinear fully linear models. In section 4 we showed how RBFs can fit in this framework, and we introduced procedures for bounding an RBF model's Hessian. In particular, these results show that the ORBIT algorithm introduced in [38] converges to first-order critical points.

The central element of an RBF is the radial function. We have illustrated the results with a few different types of radial functions. However, the results presented here are wide-reaching, requiring only the following conditions on ϕ :

1. ϕ is twice continuously differentiable on $[0, u)$, for some $u > 0$,
2. $\phi'(0) = 0$, and
3. ϕ is conditionally positive definite of order 2.

While the last condition seems to be the most restrictive, only the first condition eliminates the thin-plate spline, popular in other applications of RBFs, from our analysis. Indeed, the numerical results show that the thin-plate spline performed worst among the tested variants. We anticipate that this very general framework will be useful to researchers developing new optimization algorithms based on RBFs. Indeed, this theory extends to both the BOOSTERS algorithm [30] and ORBIT algorithm [38].

Our numerical results are aimed at illustrating the effect of using different types of radial functions ϕ in the ORBIT algorithm [38]. We saw that the cubic radial function slightly outperformed the multiquadric radial function, while the Gaussian radial function performed worse. These results are interesting because Gaussian radial basis functions are the only ones among those tested that are conditionally positive definite of order 0, requiring neither a linear nor a constant term to uniquely interpolated scattered data. Gaussian RBFs are usually used in kriging [11], which forms the basis for the global optimization methods such as [19]. We also found that the performance differences are greater when the RBF type is changed than when the maximum number of interpolation points is varied.

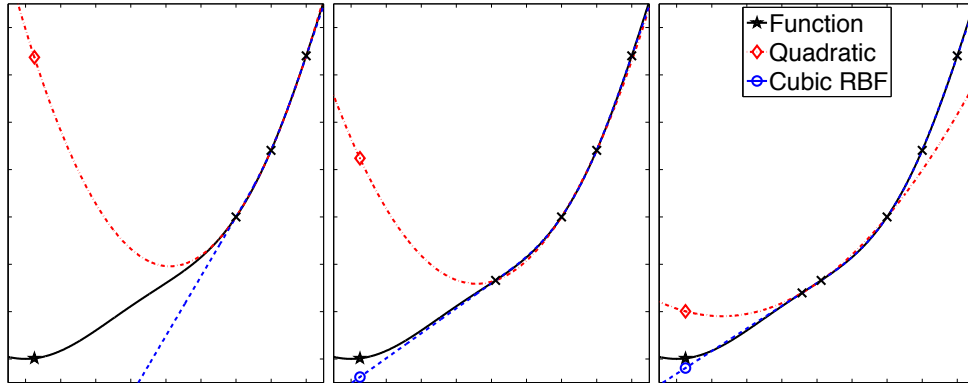


FIG. 6.1. The function $f(x) = x \sin(x\pi/4) + x^2$ approximated by a quadratic interpolating $\frac{(n+1)(n+2)}{2} = 3$ points and a cubic RBF interpolating (from left to right) 3, 4, and 5 points.

We also ran ORBIT on a computationally expensive environmental engineering problem, requiring 3 CPU-days for a single run of 80 evaluations. On this problem ORBIT quickly found a local minimum and obtained a good solution within 50 expensive evaluations.

Not surprisingly, there is no “free lunch”: while a method using RBFs outperformed methods using quadratics on the two application problems in [38], a quadratic method found the best solution on the application considered here when given a large enough budget of evaluations. Determining when to use a quadratic and when to use an RBF remains an open research problem. Our experience suggests that RBFs can be especially useful when f is nonconvex and has nontrivial higher-order derivatives.

An example of how this difference is amplified as more interpolation points are allowed is shown in Figure 6.1. As the number of points interpolated grow, the RBF model exhibits better extrapolation than does the quadratic with a fixed number of points. Similar behavior is seen even when the additional points are incorporated using a regression quadratic or a higher-order polynomial.

The present work focused primarily on the theoretical implications needed to ensure that methods using radial basis function models fit in a globally convergent trust-region framework. The results on the Blaine problem and the behavior seen in Figure 6.1 have motivated our development of global optimization methods in [36], and we intend to pursue “large-step” variants of ORBIT designed to step over computational noise.

We note that the theory presented here can be extended to models of other forms. We mention quadratics in [37], but we could also have used higher-order polynomial tails for better approximation bounds. For example, methods using a suitably conditioned quadratic tail could be expected to converge to second-order local minima. In fact, we attribute the quadratic-like convergence behavior RBF methods exhibit when at least $\frac{(n+1)(n+2)}{2}$ points are interpolated to the fact that the RBF models are *fully quadratic* with probability 1, albeit with theoretically large Taylor constants. We leave as future work the extensive numerical testing needed when many points are interpolated.

Acknowledgments. We are grateful to Amandeep Singh for providing the application problem code and to Rommel Regis and two anonymous referees for comments

on an earlier draft.

REFERENCES

- [1] C. AUDET AND J.E. DENNIS, JR., *Mesh adaptive direct search algorithms for constrained optimization*, SIAM J. Optim., 17 (2006), pp. 188–217.
- [2] A.S. BANDEIRA, K. SCHEINBERG, AND L.N. VICENTE, *Computation of sparse low degree interpolating polynomials and their application to derivative-free optimization*, Math. Program., 134 (2012), pp. 223–257.
- [3] D. BECKER, B. MINSKER, R. GREENWALD, Y. ZHANG, K. HARRE, K. YAGER, C. ZHENG, AND R. PERALTA, *Reducing long-term remedial costs by transport modeling optimization*, Ground Water, 44 (2006), pp. 864–875.
- [4] M. BJÖRKMAN AND K. HOLMSTRÖM, *Global optimization of costly nonconvex functions using radial basis functions*, Optim. Eng., 1 (2000), pp. 373–397.
- [5] M.D. BUHMANN, *Radial Basis Functions: Theory and Implementations*, Cambridge University Press, Cambridge, UK, 2003.
- [6] A.R. CONN, N.I.M. GOULD, AND P.L. TOINT, *Trust-Region Methods*, SIAM, Philadelphia, PA, 2000.
- [7] A.R. CONN, K. SCHEINBERG, AND P.L. TOINT, *Recent progress in unconstrained nonlinear optimization without derivatives*, Math. Program., 79 (1997), pp. 397–414.
- [8] A.R. CONN, K. SCHEINBERG, AND L.N. VICENTE, *Geometry of interpolation sets in derivative free optimization*, Math. Program., 111 (2008), pp. 141–172.
- [9] A.R. CONN, K. SCHEINBERG, AND L.N. VICENTE, *Global convergence of general derivative-free trust-region algorithms to first and second order critical points*, SIAM J. Optim., 20 (2009), pp. 387–415.
- [10] A.R. CONN, K. SCHEINBERG, AND L.N. VICENTE, *Introduction to Derivative-Free Optimization*, SIAM, Philadelphia, PA, 2009.
- [11] N.A. CRESSIE, *Statistics for Spatial Data*, John Wiley, New York, 1993.
- [12] A.L. CUSTÓDIO, H. ROCHA, AND L.N. VICENTE, *Incorporating minimum Frobenius norm models in direct search*, Comput. Optim. Appl., 46 (2009), pp. 265–278.
- [13] A.L. CUSTÓDIO AND L.N. VICENTE, *Using sampling and simplex derivatives in pattern search methods*, SIAM J. Optim., 18 (2007), pp. 537–555.
- [14] G. FASANO, J.L. MORALES, AND J. NOCEDAL, *On the geometry phase in model-based algorithms for derivative-free optimization*, Optim. Methods Softw., 24 (2009), pp. 145–154.
- [15] A. GRIEWANK AND A. WALTHER, *Evaluating Derivatives: Principles and Techniques of Algorithmic Differentiation, Second Ed.*, SIAM, Philadelphia, PA, 2008.
- [16] H.-M. GUTMANN, *A radial basis function method for global optimization*, J. Global Optim., 19 (2001), pp. 201–227.
- [17] N.J. HIGHAM, *The Matrix Computation Toolbox*; available online from www.ma.man.ac.uk/~higham/mctoolbox.
- [18] P.D. HOUGH, T.G. KOLDA, AND V.J. TORCZON, *Asynchronous parallel pattern search for nonlinear optimization*, SIAM J. Sci. Comput., 23 (2001), pp. 134–156.
- [19] D.R. JONES, M. SCHONLAU, AND W.J. WELCH, *Efficient global optimization of expensive black-box functions*, J. Global Optim., 13 (1998), pp. 455–492.
- [20] C.T. KELLEY, *Implicit Filtering*, SIAM, Philadelphia, PA, 2011.
- [21] T.G. KOLDA, *Revisiting asynchronous parallel pattern search for nonlinear optimization*, SIAM J. Optim., 16 (2005), pp. 563–586.
- [22] T.G. KOLDA, R.M. LEWIS, AND V.J. TORCZON, *Optimization by direct search: New perspectives on some classical and modern methods*, SIAM Rev., 45 (2003), pp. 385–482.
- [23] J.C. LAGARIAS, J.A. REEDS, M.H. WRIGHT, AND P.E. WRIGHT, *Convergence properties of the Nelder–Mead simplex algorithm in low dimensions*, SIAM J. Optim., 9 (1998), pp. 112–147.
- [24] M. MARAZZI AND J. NOCEDAL, *Wedge trust region methods for derivative free optimization*, Math. Program., 91 (2002), pp. 289–305.
- [25] J.J. MORÉ AND D.C. SORENSEN, *Computing a trust region step*, SIAM J. Sci. Stat. Comp., 4 (1983), pp. 553–572.
- [26] J.J. MORÉ AND S.M. WILD, *Benchmarking derivative-free optimization algorithms*, SIAM J. Optim., 20 (2009), pp. 172–191.
- [27] J.J. MORÉ AND S.M. WILD, *Estimating derivatives of noisy simulations*, ACM Trans. Math. Soft., 38 (2012), article 19.
- [28] T.S. MUNSON, J. SARICH, S.M. WILD, S. BENSON, AND L. CURFMAN MCINNES, *TAO 2.1 Users Manual*, Argonne Nat. Lab. Tech. Memo. ANL/MCS-TM-322, 2012; available online from

- <http://www.mcs.anl.gov/tao>.
- [29] R. OEUVRAY, *Trust-Region Methods Based on Radial Basis Functions with Application to Biomedical Imaging*, Ph.D. thesis, EPFL, Lausanne, Switzerland, 2005.
 - [30] R. OEUVRAY AND M. BIERLAIRE, *BOOSTERS: A derivative-free algorithm based on radial basis functions*, Int. J. Model. Simul., 29 (2009), pp. 26–36.
 - [31] M.J.D. POWELL, *UOBYQA: Unconstrained optimization by quadratic approximation*, Math. Program., 92 (2002), pp. 555–582.
 - [32] M.J.D. POWELL, *The NEWUOA software for unconstrained optimization without derivatives*, in Large-Scale Nonlinear Optimization, G. Di Pillo and M. Roma, eds., Springer, Berlin, 2006, pp. 255–297.
 - [33] R.G. REGIS AND C.A. SHOEMAKER, *A stochastic radial basis function method for the global optimization of expensive functions*, INFORMS J. Comput., 19 (2007), pp. 457–509.
 - [34] U.S. GEOLOGICAL SURVEY, *MODFLOW 2000*, 2000; available online from <http://water.usgs.gov/nrp/gwsoftware/modflow2000/modflow2000.html>.
 - [35] H. WENDLAND, *Scattered Data Approximation*, Cambridge Monogr. Appl. Comput. Math., Cambridge University Press, Cambridge, UK, 2005.
 - [36] S.M. WILD, *Derivative-Free Optimization Algorithms for Computationally Expensive Functions*, Ph.D. thesis, Cornell University, Ithaca, NY, 2008.
 - [37] S.M. WILD, *MNH: A derivative-free optimization algorithm using minimal norm Hessians*, in Proceedings of the Tenth Copper Mountain Conference on Iterative Methods, 2008; available online from <http://grandmaster.colorado.edu/~copper/2008/SCWinners/Wild.pdf>.
 - [38] S.M. WILD, R.G. REGIS, AND C.A. SHOEMAKER, *ORBIT: Optimization by radial basis function interpolation in trust-regions*, SIAM J. Sci. Comput., 30 (2008), pp. 3197–3219.
 - [39] Y. ZHANG, R. GREENWALD, B. MINSKER, R. PERALTA, C. ZHENG, K. HARRE, D. BECKER, L. YEH, AND K. YAGER, *Final Cost and Performance Report Application of Flow and Transport Optimization Codes to Groundwater Pump and Treat Systems*, Tech. Report TR-2238-ENV, Naval Facilities Engineering Service Center, Port Hueneme, CA, 2004.
 - [40] C. ZHENG AND P.P. WANG, *MT3DMS, a Modular Three-Dimensional Multi-species Transport Model for Simulation of Advection, Dispersion and Chemical Reactions of Contaminants in Groundwater Systems—Documentation and User’s Guide*, Contract Report SERDP-99-1, U.S. Army Engineer Research and Development Center, Vicksburg, MS, 1999.